

PATENT

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Application of: Sachin Mullick, et al.

Serial No. 10/668,467

Confirm 2942

Filed: 09/23/2003

For: Multi-Threaded Write Interface and
Methods for Increasing the Single File Read
and Write Throughput of a File Server

Group Art Unit: 2161

Examiner: Bibbee, Jared M.

Atty. Dkt. No.: 10830.0100.NPUS00

RULE 131 DECLARATION OF SORIN FAIBISH

Commissioner for Patents
PO Box 1450
Alexandria, VA 22313-1450

Sir:

1. I am a joint inventor of the subject patent application Ser. 10/668,467. I am presently employed by EMC Corporation, and I have been employed by EMC Corporation since at least 1997.

2. In the course of my work for EMC Corporation, together with Sachin Mullick, Jiannan Zheng, Xiaoye Jiang, and Peter Bixby, I was involved in the development of an invention related to a multi-threaded write interface, and the preparation of a patent application on this invention. I was responsible for reviewing and revising a first draft of the patent application in order to obtain a subsequent draft suitable for circulation among the other inventors.

3. Some time prior to Aug. 8, 2003, I received, via electronic mail, a first draft of a patent application on my multi-threaded write interface invention from my patent attorney, Richard C. Auchterlonie. Attached in Exhibit A is a true and correct copy of the transmittal E-mail message that I received with this first draft of the patent application, except that the date of the transmittal E-mail message has been redacted, where indicated by the box labeled "REDACTED".

4. On August 8, 2003, I completed a revision of the first draft of the patent application on my multi-threaded write interface invention, and I transmitted, via electronic mail, my revised version of the first draft of the patent application to my patent attorney, Richard C. Auchterlonie. Attached in Exhibit B is a true and correct copy of the transmittal E-mail and my revised version of the first draft of the patent application that was attached to this electronic mail on August 8, 2003, except that attorney-client communication has been redacted from the body of the E-mail, and attorney-client communication has been redacted from page 41 of my revised version of the first draft of the patent application, where indicated by boxes labeled "REDACTED".

5. On September 4, 2003, I received, via electronic mail, a second draft of a patent application on my multi-threaded write interface invention from my patent attorney, Richard C. Auchterlonie. Attached in Exhibit C is a true and correct copy of the transmittal E-mail that I received on September 4, 2003. I reviewed this second draft and circulated this second draft to Sachin Mullick, Jiannan Zheng, Xiaoye Jiang, and Peter Bixby, and we approved the filing of this second draft with the Patent and Trademark Office.

6. I hereby declare that all statements made of my own knowledge are true and that all statements made on information and belief are believed to be true, and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

Respectfully submitted,

Sorin Faibish / 04/09/2009
Sorin Faibish date

Linette S. Kelley

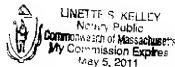


EXHIBIT A

Auchterlonie, Richard

From: Auchterlonie, Richard
Sent: REDACTED
To: staibisn@emc.com
Cc: 'Clark_William@emc.com'; 'Mazzarella_Julie@emc.com'
Subject: EMC-03-062 A New Method for Enabling Simultaneous Parallel Writes to a Single File

Sorin:

Please find attached a first draft of the specification, drawings, declaration, and assignment for your patent application. Please review and pass it along to the other inventors as appropriate. If you would like any changes, please let me know and I will send you a revised draft. Otherwise, once all of the inventors have approved of the patent application, please contact Julie Mazzarella to arrange for execution of the patent application.

Thanks,

Richard C. Auchterlonie
Howrey Simon Arnold & White, LLP
750 Bering Drive
Houston, Texas 77057
Phone: (713) 787-1698
Fax: (713) 787-1440
AuchterlonieR@howrey.com



EMCR 100 PA
Parallel Writes.DO...



EMCR 100
Drawings.pdf



EMCR100
Declaration.DOC



EMCR100
Assignment.DOC

This email message and any files transmitted with it are subject to attorney-client privilege and contains confidential information intended only for the person(s) to whom this email message is addressed. If you have received this email message in error, please notify the sender immediately by telephone or email and destroy the original message without making a copy.

EXHIBIT B

Auchterlonie, Richard

From: Sorin Faibish [sfaibish@emc.com]
Sent: Friday, August 08, 2003 4:35 PM
To: Auchterlonie, Richard
Cc: 'Clark_William@emc.com'; 'Mazzarella_Julie@emc.com'
Subject: Re: EMC-03-062 A New Method for Enabling Simultaneous Parallel Writes to a Single File



ParallelWritesPaten
tReview.zip...

Richard,

Attached please find the review of the above patent.

REDACTED

Thank you very much for your help.

"Auchterlonie, Richard" wrote:

> Sorin:
>
> Please find attached a first draft of the specification, drawings,
> declaration, and assignment for your patent application. Please review and
> pass it along to the other inventors as appropriate. If you would like any
> changes, please let me know and I will send you a revised draft. Otherwise,
> once all of the inventors have approved of the patent application, please
> contact Julie Mazzarella to arrange for execution of the patent application.
>
> Thanks,
>
> Richard C. Auchterlonie
> Howrey Simon Arnold & White, LLP
> 750 Bering Drive
> Houston, Texas 77057
> Phone: (713) 787-1698
> Fax: (713) 787-1440
> AuchterlonieR@howrey.com
>
> <<EMCR 100 PA Parallel Writes.DOC>> <<EMCR 100 Drawings.pdf>> <<EMCR100
> Declaration.DOC>> <<EMCR100 Assignment.DOC>>
>
> This email message and any files transmitted with it are subject to
> attorney-client privilege and contains confidential information intended
> only for the person(s) to whom this email message is addressed. If you have
> received this email message in error, please notify the sender immediately
> by telephone or email and destroy the original message without making a

PATENT

10830.0100.NPUS00

APPLICATION FOR UNITED STATES LETTERS PATENT

for

MULTI-THREADED WRITE INTERFACE AND
METHODS FOR INCREASING THE SINGLE FILE READ AND WRITE
THROUGHPUT OF A FILE SERVER

By

Sachin Mullick

Jiannan Zheng

Xiaoye Jiang

Sorin Faibish

Peter Bixby

Express Mail Mailing Label No. _____

BACKGROUND OF THE INVENTION

1. Field of the Invention

The present invention relates generally to file servers and data processing networks. The present invention more particularly relates to a file server permitting concurrent writes from multiple clients to the same file. The invention specifically relates to increasing the single file write throughput of such a file server.

2. Description of the Related Art

Network data storage is most economically provided by an array of low-cost disk drives integrated with a large semiconductor cache memory. A number of data mover computers are used to interface the cached disk array to the network. The data mover computers perform file locking management and mapping of the network files to logical block addresses of storage in the cached disk array, and move data between network clients and the storage in the cached disk array.

Data consistency problems may arise if multiple clients or processes have concurrent access to read-write files. Typically write synchronization and file locking have been used to ensure data consistency. For example, the data write path for a file has been serialized by holding an exclusive lock on the file for the entire duration of creating a list of data buffers to be written to disk, allocating the actual on-disk storage, and writing to storage synchronously. Unfortunately, these methods involve considerable access delays due to contention for locks not only on the files but also on the file directories and a log used when committing data to storage. In order to reduce these delays, a file server may permit asynchronous writes in accordance with version 3 of the Network File System (NFS) protocol. See, for example, Vahalia et al. U.S. Patent

1 5,893,140 issued April 6, 1999, entitled "File Server Having a File System Cache and
2 Protocol for Truly Safe Asynchronous Writes," incorporated herein by reference. More
3 recently, byte range locking to a file has been proposed in version 4 of the NFS protocol.
4 (See NFS Version 3 Protocol Specification, RFC 1813, Sun Microsystems, Inc., June
5 1995, incorporated herein by reference, and NFS Version 4 Protocol Specification, RFC
6 3530, Sun Microsystems, Inc., April 2003, incorporated herein by reference.)

7 Asynchronous writes and range locking alone will not eliminate access delays due
8 to contention during allocation and commitment of file metadata. A Unix-based file in
9 particular contains considerable metadata in the inode for the file and in indirect blocks of
10 the file. The inode, for example, contains the date of creation, date of access, file name,
11 and location of the data blocks used by the file in bitmap format. The NFS protocol
12 specifies how this metadata must be managed. In order to comply with the NFS protocol,
13 each time a write operation occurs, access to the file is not allowed until the metadata is
14 updated on disk, both for read and write operations. In a network environment, multiple
15 clients may issue simultaneous writes to the same large file such as a database, resulting
16 in considerable access delay during allocation and commitment of file metadata.

17

18 **SUMMARY OF THE INVENTION**

19 In accordance with one aspect of the present invention, there is provided a method
20 of operating a network file server for providing clients with concurrent write access to a
21 file. The method includes the network file server responding to a concurrent write
22 request from a client by obtaining a lock for the file, and then preallocating a metadata
23 block for the file, and then releasing the lock for the file, and then asynchronously writing

1 to the file, and then obtaining the lock for the file, and then committing the metadata
2 block to the file, and then releasing the lock for the file.

3 In accordance with another aspect, the invention provides a method of operating a
4 network file server for providing clients with concurrent write access to a file. The
5 method includes the network file server responding to a concurrent write request from a
6 client by preallocating a block for the file, and then asynchronously writing to the file,
7 and then committing the block to the file. The asynchronous writing to the file includes a
8 partial write to a new block that has been copied at least in part from an original block of
9 the file. The method further includes checking a partial block conflict queue for a
10 conflict with a concurrent write to the new block, and upon finding an indication of a
11 conflict with a concurrent write to the new block, waiting until resolution of the conflict
12 with the concurrent write to the new block, and then performing the partial write to the
13 new block.

14 In accordance with another aspect, the invention provides a method of operating a
15 network file server for providing clients with concurrent write access to a file. The
16 method includes the network file server responding to a concurrent write request from a
17 client by preallocating a metadata block for the file, and then asynchronously writing to
18 the file, and then committing the metadata block to the file. The method further includes
19 gathering together preallocated metadata blocks for a plurality of client write requests to
20 the file, and committing together the preallocated metadata blocks for the plurality of
21 client write requests to the file by obtaining a lock for the file, committing the gathered
22 preallocated metadata blocks for the plurality of client write requests to the file, and then
23 releasing the lock for the file.

1 In accordance with yet another aspect, the invention provides a method of
2 operating a network file server for providing clients with concurrent write access to a file.
3 The method includes the network file server responding to a concurrent write request
4 from a client by executing a write thread. Execution of the write thread includes
5 obtaining an allocation mutex for the file, and then preallocating new metadata blocks
6 that need to be allocated for writing to the file, and then releasing the allocation mutex for
7 the file, and then issuing asynchronous write requests for writing to the file, waiting for
8 callbacks indicating completion of the asynchronous write requests,
9 obtaining the allocation mutex for the file, and then committing the preallocated metadata
10 blocks, and then releasing the allocation mutex for the file.

11 In accordance with another aspect, the invention provides a network file server.
12 The network file server includes storage for storing a file and at least one processor
13 coupled to the storage for providing clients with concurrent write access to the file. The
14 network file server is programmed for responding to a concurrent write request from a
15 client by obtaining a lock for the file, and then preallocating a metadata block for the file,
16 and then releasing the lock for the file, and then asynchronously writing to the file, and
17 then obtaining the lock for the file, and then committing the metadata block to the file,
18 and then releasing the lock for the file.

19 In accordance with another aspect, the invention provides a network file server.
20 The network file server includes storage for storing a file, and at least one processor
21 coupled to the storage for providing clients with concurrent write access to the file. The
22 network file server is programmed for responding to a concurrent write request from a
23 client by preallocating a block for the file, and then asynchronously writing to the file,

1 and then committing the block to the file. The network file server includes a partial block
2 conflict queue for indicating a concurrent write to a new block that is being copied at
3 least in part from an original block of the file. The network file server is programmed for
4 responding to a client request for a partial write to the new block by checking the partial
5 block conflict queue for a conflict, and upon finding an indication of a conflict, waiting
6 until resolution of the conflict with the concurrent write to the new block of the file, and
7 then performing the partial write to the new block of the file.

8 In accordance with another aspect, the invention provides a network file server.
9 The network file server includes storage for storing a file, and at least one processor
10 coupled to the storage for providing clients with concurrent write access to the file. The
11 network file server is programmed for responding to a concurrent write request from a
12 client by preallocating a metadata block for the file, and then asynchronously writing to
13 the file; and then committing the metadata block to the file. The network file server is
14 programmed for gathering together preallocated metadata blocks for a plurality of client
15 write requests to the file, and committing together the preallocated metadata blocks for
16 the plurality of client write requests to the file by obtaining a lock for the file, committing
17 the gathered preallocated metadata blocks for the plurality of client write requests to the
18 file, and then releasing the lock for the file.

19 In accordance with yet still another aspect, the invention provides a network file
20 server. The network file server includes storage for storing a file, and at least one
21 processor coupled to the storage for providing clients with concurrent write access to the
22 file. The network file server is programmed with a write thread for responding to a
23 concurrent write request from a client by obtaining an allocation mutex for the file, and

1 then preallocating new metadata blocks that need to be allocated for writing to the file,
2 and then releasing the allocation mutex for the file, and then issuing asynchronous write
3 requests for writing to the file, waiting for callbacks indicating completion of the
4 asynchronous write requests, and then obtaining the allocation mutex for the file, and
5 then committing the preallocated metadata blocks, and then releasing the allocation
6 mutex for the file.

7 In accordance with a final aspect, the invention provides a network file server.
8 The network file server includes storage for storing a file, and at least one processor
9 coupled to the storage for providing clients with concurrent write access to the file. The
10 network file server is programmed for responding to a concurrent write request from a
11 client by preallocating a block for writing to the file, asynchronously writing to the file,
12 and then committing the preallocated block. The network file server also includes an
13 uncached write interface, a file system cache, and a cached read-write interface. The
14 uncached write interface bypasses the file system cache for sector-aligned write
15 operations, and the network file server is programmed to invalidate cache blocks in the
16 file system cache including sectors being written to by the cached read-write interface.

17

18 **BRIEF DESCRIPTION OF THE DRAWINGS**

19 Other objects and advantages of the invention will become apparent upon reading
20 the following detailed description with reference to the accompanying drawings wherein:

21 FIG. 1 is a block diagram of a data processing system including multiple clients
22 and a network file server;

FIG. 2 is a block diagram showing further details of the network file server in the data processing system of FIG. 1;

FIG. 3 is a block diagram of various read and write interfaces in a Unix-based file system layer (UxFS) in the network file server of FIG. 2;

FIG. 4 shows various file system data structures associated with a file in the network file server of FIG. 2;

FIGS. 5 and 6 comprise a flowchart of programming in the Common File System (CFS) layer in the network file server for handling a write request from a client;

FIG. 50 comprise a diagram showing multiple read or write I/Os pipelined into parallel streams in the Common File System (CFS) layer in the network file server for handling concurrent read and write requests from a client;

FIG. 5a comprise a flowchart of programming in the Common File System (CFS) layer in the network file server for handling a read request from a client;

FIG. 5b comprise a flowchart of programming in the Common File System (CFS) layer in the network file server for handling concurrent read and write requests from a client;

FIG. 7 is a flowchart of a write thread in the UxFS layer of the network file server;

FIG. 8 is a more detailed flowchart of steps in the write thread for committing preallocated metadata;

FIG. 9 is a block diagram of a partial block write during a copy-on-write operation;

FIG. 10 is a block diagram of a read-write file as maintained by the UxFS layer;

FIG. 11 is a block diagram of the read-write file of FIG. 10 after creation of a read-only version of read-write file;

FIG. 12 is a block diagram of the read-write file of FIG. 11 after a copy-on-write operation upon a direct block and two indirect blocks between the direct block and the inode of the read-write file;

FIG. 13 is a flowchart of steps in a write thread for performing the partial block write operation of FIG. 9; and

FIG. 14 shows a flowchart of steps in a write thread for allocating file blocks when writing to a file having read-only versions.

While the invention is susceptible to various modifications and alternative forms, specific embodiments thereof have been shown in the drawings and will be described in detail. It should be understood, however, that it is not intended to limit the invention to the particular forms shown, but on the contrary, the intention is to cover all modifications, equivalents, and alternatives falling within the scope of the invention as defined by the appended claims.

DESCRIPTION OF ILLUSTRATIVE EMBODIMENTS

FIG. 1 shows an Internet Protocol (IP) network 20 including a network file server 21 and multiple clients 23, 24, 25. The network file server 21, for example, has multiple data mover computers 26, 27, 28 for moving data between the IP network 20 and a cached disk array 29. The network file server 21 also has a control station 30 connected via a dedicated dual-redundant data link 31 among the data movers for configuring the data movers and the cached disk array 29.

Further details regarding the network file server 21 are found in Vahalia et al., U.S. Patent 5,893,140, incorporated herein by reference, and Xu et al., U.S. Patent 6,324,581, issued Nov. 27, 2001, incorporated herein by reference. The network file server 21 is managed as a dedicated network appliance, integrated with popular network operating systems in a way, which, other than its superior performance, is transparent to the end user. The clustering of the data movers 26, 27, 28 as a front end to the cached disk array 29 provides parallelism and scalability. Each of the data movers 26, 27, 28 is a high-end commodity computer, providing the highest performance appropriate for a data mover at the lowest cost. The data mover computers 26, 27, 28 may communicate with the other network devices using standard file access protocols such as the Network File System (NFS) or the Common Internet File System (CIFS) protocols, but the data mover computers do not necessarily employ standard operating systems. For example, the network file server 21 is programmed with a Unix-based file system that has been adapted for rapid file access and streaming of data between the cached disk array 29 and the data network 20 by any one of the data mover computers 26, 27, 28.

FIG. 2 shows software modules in the data mover 26 introduced in FIG. 1. The data mover 26 has a Network File System (NFS) module 41 for supporting communication among the clients and data movers of FIG. 1 over the IP network 20 using the NFS file access protocol, and a Common Internet File System (CIFS) module 42 for supporting communication over the IP network using the CIFS file access protocol. The NFS module 41 and the CIFS module 42 are layered over a Common File System (CFS) module 43, and the CFS module is layered over a Universal File System

(UxFS) module 44. The UxFS module supports a UNIX-based file system, and the CFS module 43 provides higher-level functions common to NFS and CIFS.

The UxFS module accesses data organized into logical volumes defined by a module 45. Each logical volume maps to contiguous logical storage addresses in the cached disk array 29. The module 45 is layered over a SCSI driver 46 and a Fibre-channel protocol (FCP) driver 47. The data mover 26 sends storage access requests through a network interface card/host bus adapter 48 using the SCSI protocol or the Fibre-Channel protocol, depending on the physical link between the data mover 26 and the cached disk array 29.

A network interface card 49 in the data mover 26 receives IP data packets from the IP network 20. A TCP/IP module 50 decodes data from the IP data packets for the TCP connection and stores the data in message buffers 53. For example, the UxFS layer 44 writes data from the message buffers 53 to a file system 54 in the cached disk array 29. The UxFS layer 44 also reads data from the file system 54 or a file system cache 51 and copies the data into the message buffers 53 for transmission to the network clients 23, 24, 25.

To maintain the file system 54 in a consistent state during concurrent writes to a file, the UxFS layer maintains file system data structures 52 in random access memory of the data mover 26. To enable recovery of the file system 54 to a consistent state after a system crash, the UxFS layer writes file metadata to a log 55 in the cached disk array during the commit of certain write operations to the file system 54.

FIG. 3 shows various read and write interfaces in the UxFS layer. These interfaces include a cached read/write interface 61 for accessing the file system cache 51, an uncached multi-threaded write interface 63, and an uncached read interface 64.

The cached read/write interface 61 permits reads and writes to the file system cache 51. If data to be accessed does not reside in the cache, it is staged from the file system 54 to the file system cache 51. Data written to the file system cache 51 from the cached read/write interface 61 is written down to the file system cache during a commit operation. The file data is written down first, followed by writing of new file metadata to the log 55 and then writing of the new metadata to the file system 54.

The uncached multi-threaded write interface 63 is used for sector-aligned writes to the file system 54. Sectors of data (e.g., 512 byte blocks) are read from the message buffers (53 in FIG. 2) and written directly to the cached disk array 29. For example, each file block is sector aligned and is 8 K bytes in length. When a sector-aligned write occurs, any cache blocks in the file system cache that include the sectors being written to are invalidated. In effect, the uncached multi-threaded write interface 63 commits file data when writing the file data to the file system 54 in storage. The uncached multi-threaded write interface 63 allows multiple concurrent writes to the same file. If a sector-aligned write changes metadata of a file such as file block allocations, then after the data of the file has been written, the new metadata is written to the log 55, and then the new metadata is written to the file system 54. The new metadata includes modifications to the file's inode, any new or modified indirect blocks, and any modified quota reservation.

The uncached read interface 64 reads sectors of data directly from the file system 54 into the message buffers (53 in FIG. 2). For example, the read request must have a

sector aligned offset and specifies a sector count for the amount of data to be read. The data can be read into multiple message buffers in one input/output operation so long as the sectors to be read are in contiguous file system blocks.

Typically, the cached read/write interface 61 is used for reading data from read-write files and from any read-only versions of the read-write files. The uncached write interface 63 is used for sector-aligned writes to read-write files. If the writes are not sector aligned, then the cached read-write interface 61 is used. The uncached read interface 64 is used for sector-aligned reads when there is no advantage to retaining the data in the file system cache 51; for example, when streaming data to a remote copy of a file.

FIG. 4 shows various file system data structures 52 associated with a file. A virtual node (VNODE)- 71 represents the file. The virtual node 71 is linked to an allocation mutex (mutually exclusive lock) 72, a partial block conflict queue 73, a partial write wait queue 74, an input-output (I/O) list 75, a staging queue 76, and preallocation block lists 77. When a file block is preallocated, it is reserved for use in the on-disk file system 54. A preallocated file block can be linked into the in-memory file block structure in the file system cache 51 as maintained by the UxFS layer 44, and later the preallocated file block can become part of the on-disk file system 54 when the preallocated file block is committed to storage. (An example of the file block structure is shown in FIG. 10.) The write threads of the uncached multi-threaded write interface (63 in FIG. 3) use the allocation mutex 72 for serializing preallocation of file metadata blocks and commitment of the preallocated metadata blocks. For a Unix-based file, the preallocated metadata blocks include new indirect blocks, which are added to the file

1 when the file is extended. As described below with reference to FIGS. 11 to 12, one or
2 more new indirect blocks may also be added to a read-write file system when processing
3 a client request to write to a direct block that is shared between the read-write file system
4 and a read-only version of the read-write file system.

5 Preallocation of the file metadata blocks under control of the allocation mutex
6 prevents multiple writers from allocating the same metadata block. The actual data write
7 is done using asynchronous callbacks within the context of the thread, and does not hold
8 any locks. Since writing to the on-disk storage takes the majority of the time, the
9 preallocation method enhances concurrency, while maintaining data integrity.

10 The preallocation method allows concurrent writes to indirect blocks within the
11 same file. Multiple writers can write to the same indirect block tree concurrently without
12 improper replication of the indirect blocks. Two different indirect blocks will not be
13 allocated for replicating the same indirect block. The write threads use the partial block
14 conflict queue 73 and the partial write wait queue 74 to avoid conflict during partial
15 block write operations, as further described below with reference to FIG. 9.

16 The I/O list 75 maps the message buffers (53 in FIG. 2) to data blocks to be
17 written. The write threads use the I/O list 75 to implement byte range locking. The data
18 blocks, for example, are 512 bytes in length providing sector-level granularity for the
19 byte range locking. Alternatively, the data block length is a multiple of the sector size.

20 In order to prevent the log (55 in FIG. 2) from becoming a bottleneck, the
21 preallocated metadata blocks for multiple write threads writing to the file at the same
22 time are committed together under the same logging lock. Committing more than one
23 allocation under one lock increases the throughput. For this purpose, a staging queue 76

1 is allocated and linked to the file virtual node 71. Preallocation block lists 77 identify the
2 respective preallocated metadata blocks for the write threads writing to the file. The
3 staging queue 76 receives pointers to the preallocation block lists 77 of the write threads
4 waiting for the allocation mutex 72 of the file for commitment of their preallocated
5 metadata blocks. For example, the staging queue 76 is a conventional circular queue, or
6 the preallocation block lists 77 are linked together into a circular list to form the staging
7 queue. There can be multiple files, and each file can have a respective staging queue
8 waiting for commitment of the file's preallocation block lists. A wait list of staging
9 queues 78 identifies the staging queues waiting for service on a first-come, first-served
10 basis.

11 From a client's view, the write operation performed by a write thread in the
12 uncached write interface is a synchronous operation. The write thread does not return an
13 acknowledgement to the client until the write data has been written down to the file
14 system in storage, and the metadata allocation has been committed to storage.

15 FIGS. 5 and 6 show programming in the Common File System (CFS) layer in the
16 network file server for handling a write request from a client. In a first step 81, if the
17 uncached multi-threaded write interface (63 in FIG. 3) is not turned on for the file
18 system, then execution branches to step 82. For example, the uncached interface can be
19 turned on or off per file system as a mount-time option. In step 82, the CFS layer obtains
20 an exclusive lock upon the file, for example by acquiring the allocation mutex (72 in FIG.
21 4) for the file. Then in step 83, the CFS layer writes a specified number of bytes from the
22 source to the file, starting at a specified byte offset, using the cached read/write interface
23 (61 in FIG. 3). The source, for example, is one or more of the message buffers (53).

1 Then in step 84, the CFS layer releases the exclusive lock upon the file, and processing of
2 the write request is finished.

3 In step 81, if the uncached multi-threaded write interface is turned on for the file
4 system, then execution continues to step 85. In step 85, if the write data specified by the
5 write request is not sector aligned (or the data size is not in multiple sectors), then
6 execution branches to step 82. Otherwise, execution continues from step 85 to step 86.

7 In step 86, the CFS layer acquires a shared lock upon the file. The shared lock
8 prevents the CFS layer from obtaining an exclusive lock upon the file for a concurrent
9 write request (e.g., in step 82). However, as described below, the shared lock upon the
10 file does not prohibit write threads in the UxFS layer from acquiring the allocation mutex
11 (72 in FIG. 4) during the preallocation of metadata blocks or during the commitment of
12 the metadata blocks.

13 In step 87, the CFS layer checks the I/O list (75 in FIG. 4) for a conflict. If there
14 is a conflicting data block on the I/O list, then execution waits until the conflicting data
15 block is flushed out of the I/O list. In certain clustered systems in which direct data
16 access to the file in the data storage is shared with other servers or clients, execution may
17 also wait in step 87 for range locks to be released by another server or client sharing
18 direct access to the file. After step 87, execution continues to step 88 in FIG. 6.

19 In step 88 of FIG. 6, the CFS layer writes the specified number of bytes from the
20 source to the file, starting at a specified sector offset, using the uncached multi-threaded
21 write interface (63 in FIG. 3). Then in step 89, the CFS layer invalidates any cached
22 entries for the file system blocks that have been written to in the file system cache (51 in
23 FIG. 3). The invalidation occurs after completion of any reads in progress to these file

1 system blocks. In step 90, the CFS layer releases the shared lock upon the file, and
2 processing of the write request is finished.

3 The parallel write architecture can be used to achieve pipelining, since the data
4 write stage does not involve any metadata interactions. Figure 50 shows how write
5 pipelining is achieved by using the preallocation, write, and commit architecture. The
6 write is divided into three steps, namely preallocation, write, and commit. The
7 preallocation stage is achieved synchronously, and an allocation mutex prevents multiple
8 preallocations from occurring simultaneously for the same file. Once the metadata
9 preallocation stage is complete, the data transfer stage can be independently handled by
10 another virtual processing unit, using Hyper threading technology (Jackson technology)
11 [[http://developer.intel.com/design/Pentium4/manuals/IntelR_PentiumR_4_Processor -](http://developer.intel.com/design/Pentium4/manuals/IntelR_PentiumR_4_Processor_Manuals.url)
12 Manuals.url]. The data transfer request can be handed over to another virtual processing
13 unit, and it does not need any interaction with the original processor that does the
14 metadata management. Separate processing units can service data read and write requests
15 generated by the master processor that handles metadata management. The write list can
16 be handed over to a separate processing unit that will then go through the write request,
17 take the data from the network packets, write it to the disk locations specified by the data
18 write request created by the master processor, and complete the data write to the disk
19 from the network packets.

20 The actual data writes happen concurrently, and can be handled by different
21 engines or virtual processing units. The final commit of the allocations are gathered
22 together and committed under the same allocation mutex. The data writes to disk are the
23 longest stage. With pipelining, the writes can be achieved continuously. This results in

1 an increase in the number of write operations that can be performed in a given time
2 period. The architecture allows the next write metadata preallocation to occur while other
3 engines are processing the data write to disk.

4 When write I/O requests arrive at the master processor or thread the request is
5 analyzed and if there are any metadata operations associated, S1 in Fig. 50, they are
6 executed by the master processor while the block write I/O is pipelined to another
7 separate processing unit. The processing unit will pipeline multiple write I/Os, S2 in Fig
8 50, and will commit all the write I/Os to the disk independently of the metadata
9 operation. At the end of the data commit process the metadata will be committed, S3 in
10 Fig 50, to the disk as well. It must be noted that the master processor is freed to perform
11 additional metadata management operations while the virtual processing unit write the
12 I/O to the disk. There could be a pool of virtual processing units that execute the write
13 tasks and they can be allocated for additional processing tasks by the master processor.
14 Some tasks are executed only by the master processor. The master processor is
15 preallocated when the data mover is rebooted. All the processing of the pipeline is based
16 on the fact that the writes are uncached and there is no contingency or locking to the files.
17 If there are any contingencies they are solved by the master processor before the writes
18 are pipelined.

19 FIG. 5a shows a flowchart of programming in the CFS layer in the network file
20 server for handling a read request from a client simultaneous with handling a write
21 request to the same file. [Please describe it in your words].

22 Fig. 5b presents the behavior of the server when there are Read –Write
23 interactions during concurrent access of multiple I/O threads to a single file.

1 More specific this describes the case when read I/O requests are sent for blocks to which
2 there are concurrent ongoing writes. There are 2 processes that are described: the
3 modifications in the write I/O flow and the read I/O flow described in the next
4 paragraphs. Fig. 5a describe the path of the reads to a file that is written to, while Fig. 5b
5 describe the concurrent writes and reads.

6 • Write request

7 Find Write Request Range -> Find partial blocks and create a partial block list ->
8 Preallocate the metadata blocks for the range of block numbers in the inode that is
9 being written -> Send the asynchronous write requests -> Wait for asynchronous write
10 request -> Get blockCommit Lock -> Commit the preallocated metadata blocks for the
11 range we have written to the inode -> release blockCommitLock -> start
12 asynchronous writes for conflict I/Os -> Find range of blocks in the buffer cache to
13 be clobbered -> Clobber the buffer Cache for the block range being committed -> if
14 active readers, mark the cache range as clobberStaleData.

15 • Read request

16 Find Read Request Range -> Is Data in Cache - Yes -> Read data from Cache -> End
17 -----No-> get blockCommit Lock -> Get the
18 committed mapping from the inode for the read request range -> release
19 blockCommitLock -> read data from disk to the buffer cache and source -> If
20 clobberStaleData flag set in this block range, clobber the buffer cache, as some write was
21 done while we were reading the data and we do not want to cache the data.

22 The above processes are based on the following assumptions:

23 -- Writes shouldn't be blocked

- 1 – Reads can read old data, new data, or a mix
- 2 – This doesn't address overlapping writes (for now).
- 3 Reads will work as implemented for normal file access non-concurrent with writes.
- 4 The *bufmap* will be examined. Valid hints will result in references on the buffer, missing
- 5 blocks will be marked as IOP (IO in Progress) and the generation count will be set to a
- 6 value associated with this read, and then a read will be started. After completing any
- 7 reads necessary the blocks that were previously marked as IOP will be cleared in one of
- 8 the following ways:
- 9 – If the slot is cleared, then it's been purged and the just completed read should not
- 10 be entered and cached.
- 11 – If it's marked as IOP then the generation count is checked:
- 12 o If the generation count is the same as we set then we cache our hint,
- 13 o otherwise we ignore this entry we can use the data to satisfy the read.
- 14 Writes will simply be allowed to proceed. At the end of the write we'll go back
- 15 in do an invalidation of the entire range of blocks written. If the slot was empty it's
- 16 ignore, if the slot had a hint it's cleared, and if the slot was IOP the IOP will be cleared
- 17 and any waiters will be awoken.
- 18 While the above only works with the assumptions listed, I believe it should be
- 19 possible to use something similar to satisfy the following assumptions:
- 20 – Writes shouldn't be blocked
- 21 – Reads can read old data or new data, but not a mix
- 22 – Overlapping writes must be controlled.

To change the processing to fulfill this set of requirements we'd first need to prevent the mixed read case. This could be done by either completing the read completely or releasing the buffers and restarting the reads when we see a conflict (case 2b above).

Further write overlaps could be prevented by adding a new type of entry in the *bufmap* which would indicate a "Write In Progress" (WIP). Each write would enter WIP in each slot in the range flushing any old entries. Subsequent writes that encountered a WIP bit within it's own range would be required to wait. Likewise attempts to read that encounter a WIP would need to block.

FIG. 7 shows a flowchart of a write thread in the UxFS layer (44 in FIG. 2). In a first step 101, the write thread gets the allocation mutex (72 in FIG. 4) for the file. Then in step 102, the write thread preallocates metadata blocks for the block range being written to the file. In step 103, the write thread releases the allocation mutex for the file.

In step 104, the write thread issues asynchronous write requests for writing to blocks of the file. For example, a list of callbacks is created. There is one callback for each asynchronous write request consisting of up to 64 K bytes of data from one or more contiguous file system blocks. An I/O list is created for each callback. The asynchronous write requests are issued asynchronously, so multiple asynchronous writes may be in progress concurrently. In step 105, the write thread waits for the asynchronous write requests to complete.

In step 106, the write thread gets the allocation mutex for the file. In step 107, the write thread commits the preallocated metadata blocks to the file system in storage. The new metadata for the file including the preallocated metadata blocks is committed by

1 being written to the log (55 in FIG. 3). File system metadata such as the file modification
2 time, however, is not committed in step 107 and is not logged. Instead, file system
3 metadata such as the file modification time is updated at a file system sync time during
4 the flushing of file system inodes. Finally, in step 108, the write thread releases the
5 allocation mutex for the file. This method of preallocating and committing metadata
6 blocks does not need any locking or metadata transactions for re-writing to allocated
7 blocks.

8 FIG. 8 is a more detailed flowchart of steps in the write thread for committing the
9 preallocated metadata. In a first step 111, if there is not a previous commit in progress,
10 then execution continues to step 112. In step 112, the thread gets the allocation mutex for
11 the file. Then in step 113, the thread writes new metadata (identified by the thread's
12 preallocation list) to the log in storage. In step 114, the thread writes the new metadata
13 (identified by the thread's preallocation list) to the file system in storage. In step 115, the
14 thread releases the allocation mutex for the file. Finally, in step 116, the thread returns an
15 acknowledgement of the write operation.

16 In step 111, if there was a previous commit in progress, then the thread inserts a
17 pointer to the threads' preallocation list onto the tail of the staging queue for the file. If
18 the staging queue was empty, then the staging queue is put on the wait list of staging
19 queues (78 in FIG. 4). The thread is suspended, waiting for a callback from servicing of
20 the staging queue. In step 118, the metadata identified by the thread's preallocation list is
21 committed when the staging queue is serviced. The staging queue is serviced by
22 obtaining the allocation mutex for the file, writing the new metadata for all of the
23 preallocation lists on the staging queue to the log in storage, then writing this new

1 metadata to the file system in storage, and then releasing the allocation mutex for the file.
2 Once servicing of the staging queue has committed the new metadata for the thread's
3 preallocation list, execution of the thread is resumed in step 116 to return an
4 acknowledgement of the write operation. After step 116, the thread is finished with the
5 write operation.

6 FIG. 9 is a block diagram of a partial block write during a copy-on-write
7 operation. Such an operation involves copying a portion of the data from an original file
8 system block 121 to a newly allocated file system block 123, and writing a new partial
9 block of data 122 to the newly allocated file system block. The portion of the data from
10 the original file system block becomes merged with the new partial block of data 122. If
11 the new partial block of data is sector aligned, then the partial block write can be
12 performed by the uncached multi-threaded write interface (63 in FIG. 3). Otherwise, if
13 the new partial block of data were not sector aligned, then the partial block write would
14 be performed by the cached read/write interface (61 in FIG. 3).

15 The copy-on-write operation may frequently occur in a file system including one
16 or more read-only file versions of a read-write file. Such a file system is described in
17 Chutani, Sailesh, et al., "The Episode File System," Carnegie Mellon University IT
18 Center, Pittsburgh, PA, June 1991, incorporated herein by reference. Each read-only file
19 version is a snapshot of the read-write file at a respective point in time. Read-only file
20 versions can be used for on-line data backup and data mining tasks.

21 In a copy-on-write file versioning method, the read-only file version initially
22 includes only a copy of the inode of the original file. Therefore the read-only file version
23 initially shares all of the data blocks as well as any indirect blocks of the original file.

1 When the original file is modified, new blocks are allocated and linked to the original file
2 inode to save the new data, and the original data blocks are retained and linked to the
3 inode of the read-only file version. The result is that disk space is saved by only saving
4 the difference between two consecutive versions. This process is shown in FIGS. 9, 10,
5 and 11.

6 FIG. 10 shows a read-write file as maintained by the UxFS layer. The file has a
7 hierarchical organization, depicted as an inverted tree. The file includes a read-write
8 inode 131, a direct block 132 and an indirect block 133 linked to the read-write inode, a
9 direct block 134 and an indirect block 135 linked to the indirect block 133, and direct
10 blocks 136 and 137 linked to the indirect block 135.

11 When a read-only version of a read-write file is created, a new inode for the read-
12 only version is allocated. The read-write file inode and file handle remain the same.
13 After allocation of the new inode, the read-write file is locked and the new inode is
14 populated from the contents of the read-write file inode. Then the read-write file inode
15 itself is modified, the transaction is committed, and the lock on the read-write file is
16 released.

17 The allocation of blocks during the copy-on-write to the read-write file raises the
18 possibility of the supply of free storage being used up after writing to a small fraction of
19 the blocks of the read-write file. To eliminate this possibility, the read-write file can be
20 provided with a “persistent reservation” mechanism so that the creation of a read-only
21 version will fail unless there can be reserved a number of free storage blocks equal to the
22 number of blocks that become shared between the read-only version and the read-write
23 file. The number of reserved blocks can be maintained as an attribute of the file. The

1 number of reserved blocks for a read-only file can be incremented as blocks become
2 shared with a read-only version, and decremented as blocks are allocated during the
3 writes to the read-write file.

4 FIG. 11 shows the read-write file of FIG. 10 after creation of a read-only version
5 of the read-write file. The read-only inode 138 is a copy of the read-write inode 131.
6 The read-write inode 131 has been modified to indicate that the direct block 132 and the
7 indirect block 133 are shared with a read-only version. For example, in the read-write
8 inode 131, the most significant bit in each of the pointers to direct block 132 and the
9 indirect block 133 have been set to indicate that the pointers point to blocks that are
10 shared with the read-write file. (The links represented by such pointers to shared blocks
11 are indicated by dotted lines in FIGS. 11 and 12.) Also, by inheritance, any and all of the
12 descendants of a shared block are also shared blocks. Routines in the UxFS layer that
13 use the pointers to locate the pointed-to file system blocks simply mask out the most
14 significant to determine the block addresses.

15 In general, for the case in which there are multiple versions of a file sharing file
16 blocks, when a file block is shared, it is desirable to designate the oldest version sharing
17 the block to be the owner of the block, and any other files to be non-owners of the block.
18 A pointer in a non-shared block pointing to a shared block will have its most significant
19 bit set if the block is not owned by the owner of the non-shared block, and will have its
20 most significant bit clear if the block is owned by the owner of the non-shared block.

21 When writing to a specified sector of a file, a search of the file block hierarchy is
22 done starting with the read-write inode, in order to find the file block containing the
23 specified sector. Upon finding a pointer indicating that the pointed-to block is shared, the

1 pointed-to block and its descendants are noted as “copy on write” blocks. If the specified
2 sector is found in a “copy on write” block, then a new file block is allocated.

3 In practice, multiple write threads are executed concurrently, so that more than
4 one concurrent write thread could determine a need to preallocate the same new file
5 block. The allocation mutex is used to serialize the allocation process so more than one
6 preallocation of a new file block does not occur. For example, once the write thread has
7 obtained the allocation mutex, the write thread then determines whether a new block is
8 needed, and if so, then the write thread preallocates the new block. The write thread may
9 obtain the allocation mutex, allocate multiple new blocks in this fashion, and then release
10 the allocation mutex. For example, to write to a direct block of a file, when the write
11 thread finds a shared block on the path in the file hierarchy down to the direct block of
12 the file, the write thread obtains the allocation mutex, and then allocates all the shared
13 blocks that it then finds down the path in the file hierarchy down to and including the
14 direct block, and then release the allocation mutex.

15 Once a new file block has been allocated, a partial block write to the new file
16 block is performed, unless the write operation writes new data to the entire block. The
17 new file block is the same type (direct or indirect) as the original “copy on write” file
18 block containing the specified sector. If the write operation writes new data to the entire
19 new file block, then no copy need be done and the new data is simply written into the
20 newly allocated block. (A partial write could be performed when the write operation
21 writes new data to the entire block, although this would not provide the best
22 performance.)

1 If the read-write inode or a block owned by the read-write file was a parent of the
2 original “copy on write” block, then the new file block becomes a child of the read-write
3 inode or the block owned by the read-write file. Otherwise, the new file block becomes
4 the child of a newly allocated indirect block. In particular, copies are made of all of the
5 “copy on write” indirect blocks that are descendants of the read-write inode and are also
6 predecessors of the original “copy on write” file block.

7 For example, assume that a write request specifies a sector found to be in the
8 direct block 137 of FIG. 11. Upon searching down the hierarchy from the read-write
9 inode 131, it is noted that indirect blocks 133 and 135 and the direct block 137 are “copy
10 on write” blocks. As shown in FIG. 12, new indirect blocks 139 and 140 and a new
11 direct block 141 have been allocated. The new direct block 141 is a copy of the original
12 direct block 136 except that it includes the new data of the write operation. The new
13 indirect block 140 is a copy of the original indirect block 135 except it has a new pointer
14 pointing to the new direct block 141 instead of the original direct block 137. The new
15 indirect block 139 is a copy of the original indirect block 133 except it has a new pointer
16 pointing to the new indirect block 140 instead of the original indirect block 135. Also,
17 the read-write inode 131 has been modified to replace the pointer to the original indirect
18 block 133 with a pointer to the new indirect block 139.

19 In some instances, a write to the read-write file will require the allocation of a
20 new direct block without any copying from an original direct block. This occurs when
21 there is a full block write, a partial block write to a hole in the file, or a partial block write
22 to an extended portion of a file. When there is a partial block write to a hole in the file or
23 a partial block write to the extended portion of a file, the partial block of new data is

1 written to the newly allocated direct block, and the remaining portion of the newly
2 allocated direct block is filled in with zero data.

3 It is possible that the UxFS layer will receive multiple concurrent writes that all
4 require new data to be written to the same newly allocated block. These multiple
5 concurrent writes need to be synchronized so that only one new block will be allocated
6 and the later one of the threads will not read old data from the original block and copy the
7 old data onto the new data from an earlier one of the threads. The UxFS layer detects the
8 first such write request and puts a corresponding entry into the partial block conflict
9 queue (73 in FIG. 4). The UxFS layer detects the second such write request, determines
10 that it is conflicting upon inspection of the partial block conflict queue, places an entry to
11 the second such write request in the partial write wait queue (74 in FIG. 4), and suspends
12 the write thread for the second such write request until the conflict is resolved.

13 FIG. 13 is a flowchart of steps in a write thread for performing the partial block
14 write operation of FIG. 9. In a first step 151 of FIG. 13, if the newly allocated file system
15 block (124 in FIG. 9) is not on the partial block conflict queue (73 in FIG. 4), then
16 execution branches to step 152. In step 152, the partial block write thread puts the new
17 block on the partial block conflict queue. In step 153, the partial block write thread
18 copies data that will not be overwritten by the partial block write, the data being copied
19 from the original file system block to the new file system block. In step 154,
20 asynchronous write operations are performed to write the new partial block of data to the
21 new block. In step 155, the partial block write thread gets the allocation mutex for the
22 file, commits the preallocated metadata (or the preallocated metadata is gathered and
23 committed upon servicing of the staging queue if a previous commit is in progress),

1 removes the new block from the partial block conflict queue, issues asynchronous writes
2 for any corresponding blocks on the partial write wait queue, and releases the allocation
3 mutex.

4 In step 151, if the newly allocated file system block was on the partial block
5 conflict queue, then execution continues to step 156. In step 156, the partial block write
6 thread puts a write callback on the partial write wait queue for the file. Then execution is
7 suspended until the callback occurs (from the completion of the asynchronous writes
8 issued in step 155). Upon resuming, in step 157, the partial block write thread gets the
9 allocation mutex for the file, commits the preallocated metadata (or the preallocated
10 metadata is gathered and committed upon servicing of the staging queue if a previous
11 commit is in progress), and releases the allocation mutex.

12 FIG. 14 shows steps in a write thread for allocating file blocks when writing to a
13 file having read-only versions. In a first step 161, if the file block being written to is not
14 shared with a read-only version, then execution branches to step 162 to write directly to
15 the block without any transaction. In other words, there is no need for allocating any
16 additional blocks.

17 In step 161, if the file block being written to is shared with a read-only version,
18 then execution continues to step 163. In step 163, if the file block being written to is an
19 indirect block, then execution branches to step 164. In step 164, a new indirect block is
20 allocated, the original indirect block content is copied to the new indirect block, and the
21 new metadata is written to the new indirect block synchronously. If the block's parent is
22 an indirect block shared with a read-only version, then a new indirect block is allocated
23 for copy-on-write of the new block pointer. Any other valid block pointers in this new

1 indirect block point to shared blocks, and therefore the most significant bit in each of
2 these other valid block pointers should be set (as indicated by the dotted line between the
3 indirect blocks 136 and 140 in FIG. 12). For example, just after the original indirect
4 block content is copied to the new indirect block, the most significant bit is set in all valid
5 block pointers in the new indirect block. As described above with respect to FIG. 12, this
6 copy-on-write may require one or more additional indirect blocks to be allocated (such as
7 indirect block 139 in FIG. 12). For example, the tree of a UxFS file may include up to
8 three levels of indirect blocks. All of the file blocks that need to be allocated can be
9 predetermined so that the allocation mutex for the file can be obtained, all of the new
10 blocks that are needed can be allocated together, and then the allocation mutex for the file
11 can be released.

12 In step 163, if the file block being written to is not an indirect block, then
13 execution continues to step 165. This is the case in which the file block being written to
14 is a direct block. In step 165, if the write to the file block is not a partial write, then
15 execution branches to step 166. In step 166, a new direct block is allocated and the block
16 of new data is written directly to the new direct block. If the original block's parent is an
17 indirect block that is shared with a read-only version, then a new indirect block is
18 allocated for copy-on-write of the new block pointer. As described above with respect to
19 FIG. 12, this copy-on-write may require one or more additional indirect blocks to be
20 allocated.

21 In step 167, for the case of a partial write, execution continues from step 156 to
22 step 167 to use the partial write technique as described above with respect to FIG. 9 and
23 FIG. 13.

1 Various parts of the programming for handling a write thread the UxFS layer have
2 been described above with reference to FIGS. 7 to 14. Following is a listing of the steps
3 in the preferred implementation of this programming.

4 1. The write thread receives a write request specifying the source and
5 destination of the data to be written. The source is specified in terms of message buffers
6 and the message buffer header size. The destination is specified in terms of an offset and
7 number of bytes to be written.

8 2. The write thread calculates the starting and ending logical block number,
9 total block count, and determines whether the starting and ending blocks are partial
10 blocks.

11 3. The write thread gets the allocation mutex for the file.

12 4. The write thread searches the file tree along a path from the file inode to
13 the destination file blocks to determine whether there are any shared blocks along this
14 path. For each such shared block, a new direct or indirect block is allocated
15 synchronously, as described above with reference to FIGS. 11, 12, and 14.

16 5. The write thread identifies partial blocks of write data using the starting
17 physical block number and the number of blocks to be written. Only the starting and
18 ending block to be written can be partial. Also, if some other thread got to these blocks
19 first, the block mapping may already exist and the “copy-on-write” will be done by the
20 prior thread. The partial block conflict queue is checked to determine whether such an
21 allocation and “copy-on-write” is being done by a prior thread. If so, the block write of
22 the present thread is added to the partial write wait queue, as described above with
23 reference to FIG. 13.

- 1 6. The write thread preallocates the metadata blocks.
- 2 7. The write thread releases the allocation mutex.
- 3 8. The write threads determine the state of the block write. The block write
- 4 can be in one of three states, namely:
- 5 1. Partial, in-progress writes. These are writes to blocks that are on the
- 6 conflict list. This write is deferred. The information to write out these
- 7 blocks is added to the partial write wait queue.
- 8 2. Whole Block Writes.
- 9 3. Partial, not-in-progress writes. These are partial writes to newly allocated
- 10 blocks, and are the first write to these blocks.
- 11 9. The I/O list is split apart if there are any non-contiguous areas to be
- 12 written.
- 13 10. Asynchronous write requests are issued for blocks in state 2 (full block
- 14 writes).
- 15 11. Synchronous read requests are issued for blocks in state 3 (Partial not-in-
- 16 progress writes).
- 17 12. Asynchronous write requests are issued for blocks in state 3.
- 18 13. The write thread waits for all writes to complete, including the ones in
- 19 state 1. The write thread waits for all asynchronous write callbacks. The asynchronous
- 20 writes for blocks in state 1 are actually issued by other threads.
- 21 14. The write thread gets the allocation mutex.

1 15. The write thread commits the preallocated metadata. The allocation lists
2 being committed are gathered together if a previous commit is in progress, and are
3 written out under the same logging lock as described above with reference to FIG. 8.

4 16. The write thread removes any blocks that the write thread had added to
5 partial block conflict queue, and issues asynchronous writes for corresponding blocks on
6 the partial write wait queue.

7 17. The write thread releases the allocation mutex. The write thread has
8 completed the write operation.

9

10 In view of the above, there have been described a multi-threaded write interface
11 for increasing the single file write throughput of a file server. The write interface allows
12 multiple concurrent writes to the same file and handles metadata updates using sector
13 level locking. The write interface provides permission management to access the data
14 blocks of the file in parallel, ensures correct use and update of indirect blocks in the tree
15 of the file, preallocates file blocks when the file is extended, and solves access conflicts
16 for simultaneous writes to the same block. The write interface preallocates file metadata
17 to prevent multiple writers from allocating the same block. For example, a write
18 operation includes obtaining a per file allocation mutex (mutually exclusive lock),
19 reserving a metadata block, releasing the allocation mutex, issuing an asynchronous write
20 request for writing to the file, waiting for the asynchronous write request to complete,
21 obtaining the allocation mutex, committing the preallocated metadata block, and
22 releasing the allocation mutex. Since no locks are held during the writing of data to the

1 on-disk storage and this data write takes the majority of the time, the method enhances
2 concurrency while maintaining data integrity.

3

1 What is claimed is:

2

3 1. A method of operating a network file server for providing clients with concurrent
4 write access to a file, the method comprising the network file server responding to a
5 concurrent write request from a client by:

6 (a) obtaining a lock for the file; and then

7 (b) preallocating a metadata block for the file; and then

8 (c) releasing the lock for the file; and then

9 (d) asynchronously writing to the file; and then

10 (e) obtaining the lock for the file; and then

11 (f) committing the metadata block to the file; and then

12 (g) releasing the lock for the file.

13

14 2. The method as claimed in claim 1, wherein the file includes a hierarchy of blocks
15 including an inode block of metadata, direct blocks of file data, and indirect blocks of
16 metadata, and wherein the metadata block for the file is an indirect block of metadata.

17

18 3. The method as claimed in claim 2, which includes copying data from an original
19 indirect block of the file to the metadata block for the file, the original indirect block of
20 the file having been shared between the file and a read-only version of the file.

21

22 4. The method as claimed in claim 1, which includes concurrent writing for more
23 than one client to the metadata block for the file.

1

2 5. The method as claimed in claim 2, wherein the asynchronous writing to the file
3 includes a partial write to a new block that has been copied at least in part from an
4 original block of the file, and wherein the method includes checking a partial block
5 conflict queue for a conflict with a concurrent write to the new block, and upon failing to
6 find an indication of a conflict with a concurrent write to the new block, preallocating the
7 new block, copying at least a portion of the original block of the file to the new block,
8 and performing the partial write to the new block.

9

10 6. The method as claimed in claim 2, wherein the asynchronous writing to the file
11 includes a partial write to a new block that has been copied at least in part from an
12 original block of the file, and wherein the method includes checking a partial block
13 conflict queue for a conflict with a concurrent write to the new block, and upon finding
14 an indication of a conflict with a concurrent write to the new block, waiting until
15 resolution of the conflict with the concurrent write to the new block, and then performing
16 the partial write to the new block.

17

18 7. The method as claimed in claim 6, which includes placing a request for the partial
19 write in a partial write wait queue upon finding an indication of a conflict with a
20 concurrent write to the new block, and performing the partial write upon servicing the
21 partial write wait queue.

22

1 8. The method as claimed in claim 1, wherein the asynchronously writing to the file
2 includes checking an input-output list for a conflicting prior concurrent access to the file,
3 and upon finding a conflicting prior concurrent access to the file, suspending the
4 asynchronous writing to the file until the conflicting prior concurrent access to the file is
5 no longer conflicting.

6

7 9. The method as claimed in claim 8, wherein the suspending of the asynchronous
8 writing to the file until the conflicting prior concurrent access is no longer conflicting
9 provides a sector-level granularity of byte range locking for concurrent write access to
10 the file.

11

12 10. The method as claimed in claim 1, wherein the metadata block for the file is
13 committed by writing the metadata block to a log in storage of the network file server.

14

15 11. The method as claimed in claim 1, which includes gathering together preallocated
16 metadata blocks for a plurality of client write requests to the file, and committing
17 together the preallocated metadata blocks for the plurality of client write requests to the
18 file by obtaining the lock for the file, committing the gathered preallocated metadata
19 blocks for the plurality of client write requests to the file, and then releasing the lock for
20 the file.

21

22 12. The method as claimed in claim 1, which includes checking whether a previous
23 commit is in progress after asynchronously writing to the file and before obtaining the

1 lock for the file for committing the metadata block to the file, and upon finding that a
2 previous commit is in progress, placing a request for committing the metadata block to
3 the file on a staging queue for the file.

4
5 13. A method of operating a network file server for providing clients with concurrent
6 write access to a file, the method comprising the network file server responding to a
7 concurrent write request from a client by:

- 8 (a) preallocating a block for the file; and then
- 9 (b) asynchronously writing to the file; and then
- 10 (c) committing the block to the file;

11 wherein the asynchronous writing to the file includes a partial write to a new
12 block that has been copied at least in part from an original block of the file, and wherein
13 the method includes checking a partial block conflict queue for a conflict with a
14 concurrent write to the new block, and upon finding an indication of a conflict with a
15 concurrent write to the new block, waiting until resolution of the conflict with the
16 concurrent write to the new block, and then performing the partial write to the new block.

17
18 14. The method as claimed in claim 13, wherein the method includes placing a
19 request for the partial write in a partial write wait queue upon finding an indication of a
20 conflict with a concurrent write to the new block, and performing the partial write upon
21 servicing the partial write wait queue.

1 15. A method of operating a network file server for providing clients with concurrent
2 write access to a file, the method comprising the network file server responding to a
3 concurrent write request from a client by:

4 (a) preallocating a metadata block for the file; and then

5 (b) asynchronously writing to the file; and then

6 (c) committing the metadata block to the file;

7 wherein the method includes gathering together preallocated metadata blocks for
8 a plurality of client write requests to the file, and committing together the preallocated
9 metadata blocks for the plurality of client write requests to the file by obtaining a lock for
10 the file, committing the gathered preallocated metadata blocks for the plurality of client
11 write requests to the file, and then releasing the lock for the file.

12
13 16. The method as claimed in claim 15, which includes checking whether a previous
14 commit is in progress after asynchronously writing to the file and before obtaining the
15 lock for the file for committing the block to the file, and upon finding that a previous
16 commit is in progress, placing a request for committing the metadata block to the file on
17 a staging queue for the file.

18
19 17. A method of operating a network file server for providing clients with concurrent
20 write access to a file, the method comprising the network file server responding to a
21 concurrent write request from a client by executing a write thread, execution of the write
22 thread including:

23 (a) obtaining an allocation mutex for the file; and then

1 (b) preallocating new metadata blocks that need to be allocated for writing to the
2 file; and then

3 (c) releasing the allocation mutex for the file; and then

4 (d) issuing asynchronous write requests for writing to the file;

5 (e) waiting for callbacks indicating completion of the asynchronous write
6 requests; and then

7 (f) obtaining the allocation mutex for the file; and then

8 (g) committing the preallocated metadata blocks; and then

9 (h) releasing the allocation mutex for the file.

10

11 18. A network file server comprising storage for storing a file, and at least one
12 processor coupled to the storage for providing clients with concurrent write access to the
13 file, wherein the network file server is programmed for responding to a concurrent write
14 request from a client by:

15 (a) obtaining a lock for the file; and then

16 (b) preallocating a metadata block for the file; and then

17 (c) releasing the lock for the file; and then

18 (d) asynchronously writing to the file; and then

19 (e) obtaining the lock for the file; and then

20 (f) committing the metadata block to the file; and then

21 (g) releasing the lock for the file.

22

1 19. The network file server as claimed in claim 18, wherein the file includes a
2 hierarchy of blocks including an inode block of metadata, direct blocks of file data, and
3 indirect blocks of metadata, and wherein the metadata block for the file is an indirect
4 block of metadata.

5
6
7
8
9
10
11 REDACTED
12
13
14
15
16
17
18
19

20 20. The network file server as claimed in claim 19, which is programmed for copying
21 data from an original indirect block of the file to the metadata block for the file, the
22 original indirect block of the file having been shared between the file and a read-only
23 version of the file.

1

2 21. The network file server as claimed in claim 18, which is programmed for
3 concurrent writing for more than one client to the metadata block for the file.

4

5 22. The network file server as claimed in claim 18, which includes a partial block
6 conflict queue for indicating a concurrent write to a new block that is being copied at
7 least in part from an original block of the file, and wherein the network file server is
8 programmed to respond to a client request for a partial write to the new block by
9 checking the partial block conflict queue for a conflict, and upon failing to find an
10 indication of a conflict, preallocating the new block, copying at least a portion of the
11 original block of the file to the new block, and performing a partial write to the new
12 block.

13

14 23. The network file server as claimed in claim 18, which includes a partial block
15 conflict queue for indicating a concurrent write to a new block that is being copied at
16 least in part from an original block of the file, and wherein the network file server is
17 programmed to respond to a client request for a partial write to the new block by
18 checking the partial block conflict queue for a conflict, and upon finding an indication of
19 a conflict, waiting until resolution of the conflict with the concurrent write to the new
20 block, and then performing the partial write to the new block.

21

22 24. The network file server as claimed in claim 23, which includes a partial write wait
23 queue, and wherein the network file server is programmed for placing a request for the

1 partial write in the partial write wait queue upon finding an indication of a conflict, and
2 performing the partial write upon servicing the partial write wait queue.

3
4 25. The network file server as claimed in claim 18, which is programmed for
5 maintaining an input-output list of concurrent reads and writes to the file, and when
6 asynchronously writing to the file, for checking the input-output list for a conflicting
7 prior concurrent access, read or write, to the file, and upon finding a conflicting prior
8 concurrent access to the file, suspending the asynchronous writing to the file until the
9 conflicting prior concurrent access to the file is no longer conflicting. The block could be
10 read or write but must specify read explicitly and the claim should cover read access by
11 removing the file lock (mutex) and change it to a range lock.

12
13 26. The network file server as claimed in claim 25, wherein the suspending of the
14 asynchronous writing to the file until the conflicting prior concurrent access is no longer
15 conflicting provides a sector-level granularity of byte range locking for concurrent write
16 access to the file, as well as concurrent reads and writes. [Probably would be better to add
17 a new claim.]

18
19 27. The network file server as claimed in claim 18, which is programmed for
20 committing the metadata block for the file by writing the metadata block to a log in the
21 storage.

1 28. The network file server as claimed in claim 18, which is programmed for
2 gathering together preallocated metadata blocks for a plurality of client requests for write
3 access to the file, and committing together the preallocated metadata blocks for the
4 plurality of client requests for access to the file by obtaining the lock for the file,
5 committing the gathered preallocated metadata blocks for the plurality of client requests
6 for write access to the file, and then releasing the lock for the file.

7

8 29. The network file server as claimed in claim 18, which includes a staging queue
9 for the file, and which is programmed for checking whether a previous commit is in
10 progress after asynchronously writing to the file and before obtaining the lock for the file
11 for committing the metadata block to the file, and upon finding that a previous commit is
12 in progress, placing a request for committing the metadata block to the file on the staging
13 queue for the file.

14

15 30. A network file server comprising storage for storing a file, and at least one
16 processor coupled to the storage for providing clients with concurrent write access to the
17 file, wherein the network file server is programmed for responding to a concurrent write
18 request from a client by:

19 (a) preallocating a block for the file; and then

20 (b) asynchronously writing to the file; and then

21 (c) committing the block to the file;

22 wherein the network file server includes a partial block conflict queue for indicating a
23 concurrent write to a new block that is being copied at least in part from an original block

1 of the file, and wherein the network file server is programmed for responding to a client
2 request for a partial write to the new block by checking the partial block conflict queue
3 for a conflict, and upon finding an indication of a conflict, waiting until resolution of the
4 conflict with the concurrent write to the new block of the file, and then performing the
5 partial write to the new block of the file.

6

7 31. The network file server as claimed in claim 30, which includes a partial write wait
8 queue, and wherein the network file server is programmed for placing a request for the
9 partial write in the partial write wait queue upon finding an indication of a conflict, and
10 performing the partial write upon servicing the partial write wait queue.

11

12 32. A network file server comprising storage for storing a file, and at least one
13 processor coupled to the storage for providing clients with concurrent write access to the
14 file, wherein the network file server is programmed for responding to a concurrent write
15 request from a client by:

16 (a) preallocating a metadata block for the file; and then

17 (b) asynchronously writing to the file; and then

18 (c) committing the metadata block to the file;

19 wherein the network file server is programmed for gathering together preallocated
20 metadata blocks for a plurality of client write requests to the file, and committing
21 together the preallocated metadata blocks for the plurality of client write requests to the
22 file by obtaining a lock for the file, committing the gathered preallocated metadata blocks
23 for the plurality of client write requests to the file, and then releasing the lock for the file.

1
2 33. The network file server as claimed in claim 32, which is programmed for
3 checking whether a previous commit is in progress after asynchronously writing to the
4 file and before obtaining the lock for the file for committing the metadata block to the
5 file, and upon finding that a previous commit is in progress, placing a request for
6 committing the metadata block to the file on a staging queue for the file.

7
8 34. A network file server comprising storage for storing a file, and at least one
9 processor coupled to the storage for providing clients with concurrent write access to the
10 file, wherein the network file server is programmed with a write thread for responding to
11 a concurrent write request from a client by:

12 (a) obtaining an allocation mutex for the file; and then

13 (b) preallocating new metadata blocks that need to be allocated for writing to the
14 file; and then

15 (c) releasing the allocation mutex for the file; and then

16 (d) issuing asynchronous write requests for writing to the file;

17 (e) waiting for callbacks indicating completion of the asynchronous write
18 requests; and then

19 (f) obtaining the allocation mutex for the file; and then

20 (g) committing the preallocated metadata blocks; and then

21 (h) releasing the allocation mutex for the file.
22

1 35. The network file server as claimed in claim 34, which includes an uncached write
2 interface, a file system cache and a cached read-write interface, and wherein the
3 uncached write interface bypasses the file system cache for sector-aligned write
4 operations.

5
6 36. The network file server as claimed in claim 35, wherein the network file server is
7 programmed to invalidate cache blocks in the file system cache including sectors being
8 written to by the cached read-write interface.

9
10 37. A network file server comprising storage for storing a file, and at least one
11 processor coupled to the storage for providing clients with concurrent write access to the
12 file, wherein the network file server is programmed for responding to a concurrent write
13 request from a client by:

14 (a) preallocating a block for writing to the file;

15 (b) asynchronously writing to the file; and then

16 (c) committing the preallocated block;

17 wherein the network file server also includes an uncached write interface, a file system
18 cache, and a cached read-write interface, wherein the uncached write interface bypasses
19 the file system cache for sector-aligned write operations, and the network file server is
20 programmed to invalidate cache blocks in the file system cache including sectors being
21 written to by the cached read-write interface.

22 [Please add a new set of claims that refer explicitly to files that have read-only copies.]

ABSTRACT

A write interface in a file server provides permission management for concurrent access to data blocks of a file, ensures correct use and update of indirect blocks in a tree of the file, preallocates file blocks when the file is extended, and solves access conflicts for simultaneous writes to the same block. For example, a write operation includes obtaining a per file allocation mutex (mutually exclusive lock), preallocating a metadata block, releasing the allocation mutex, issuing an asynchronous write request for writing to the file, waiting for the asynchronous write request to complete, obtaining the allocation mutex, committing the preallocated metadata block, and releasing the allocation mutex. Since no locks are held during the writing of data to the on-disk storage and this data write takes the majority of the time, the method enhances concurrency while maintaining data integrity. [Please add a sentence mentioning that the read are also faster as they do not have to wait for the file lock during a read. Also mention tha this is critical for database applications.]

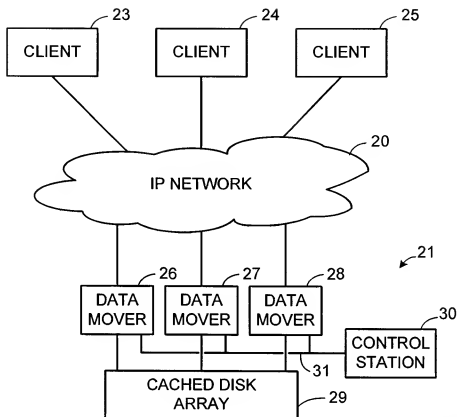


FIG. 1

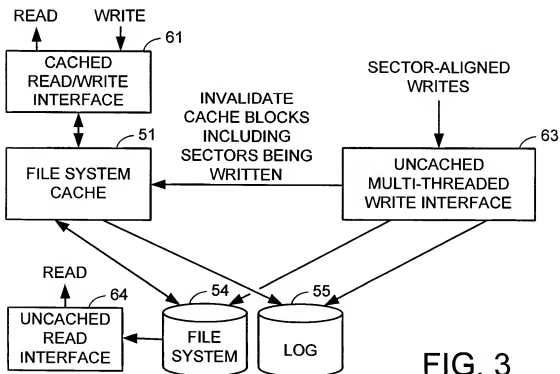


FIG. 3

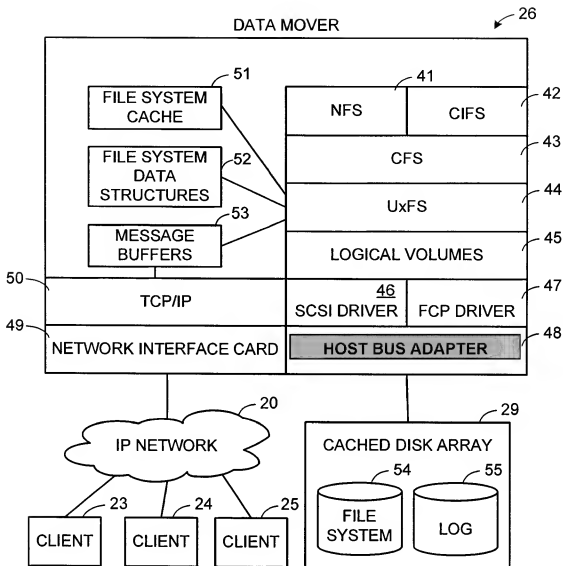


FIG. 2

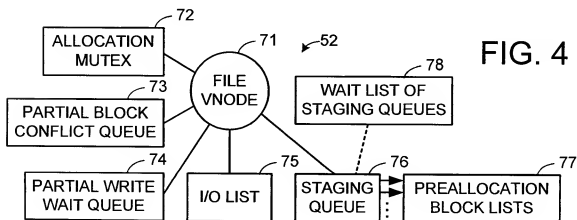


FIG. 4

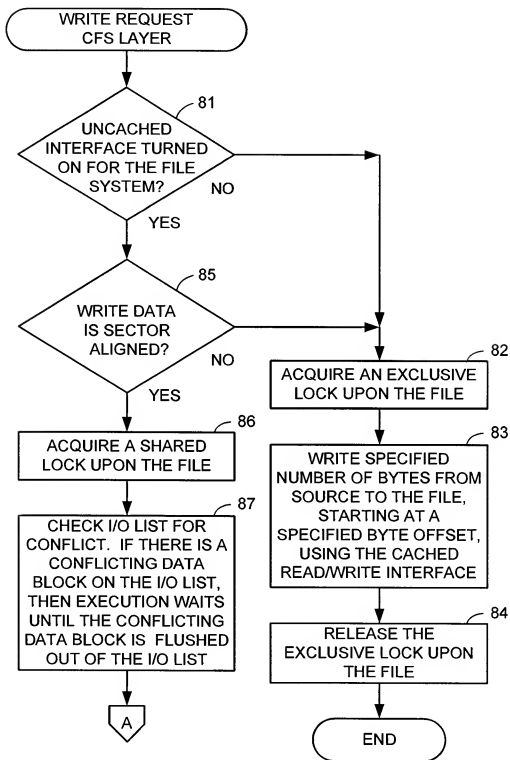


FIG. 5

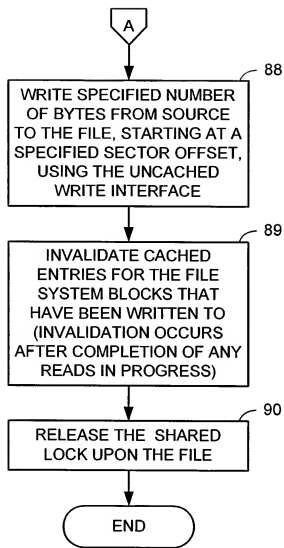


FIG. 6

NEW

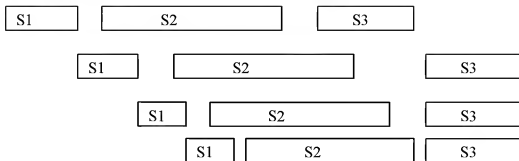


Figure 50 Write Request Pipelining: Writes are divided into three stages

S1: preallocation of metadata

S2: data writes

S3: Commit of metadata

Fig. 50

NEW

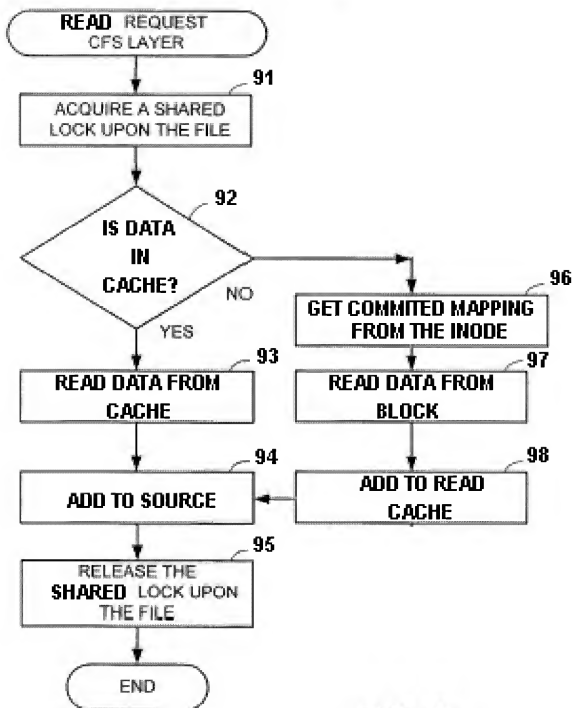


FIG. 5a

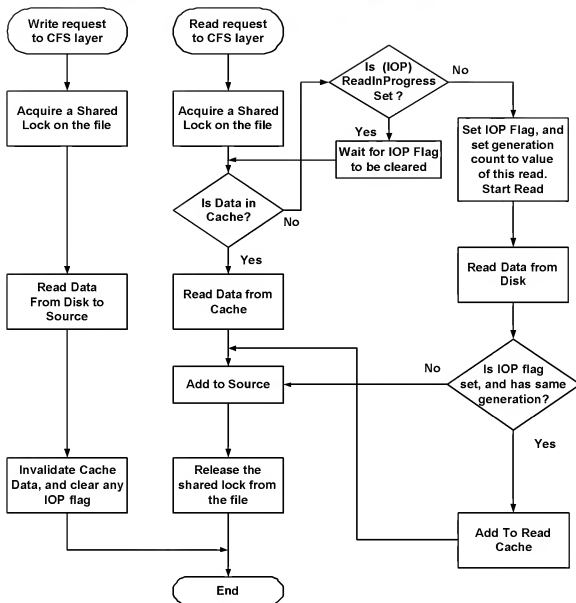


Fig 5b

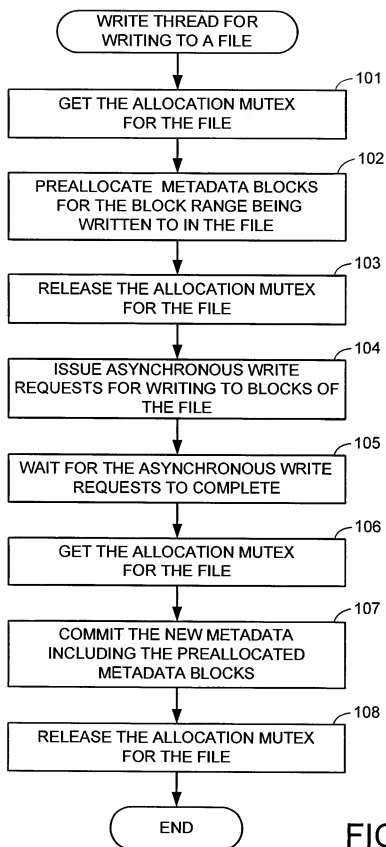


FIG. 7

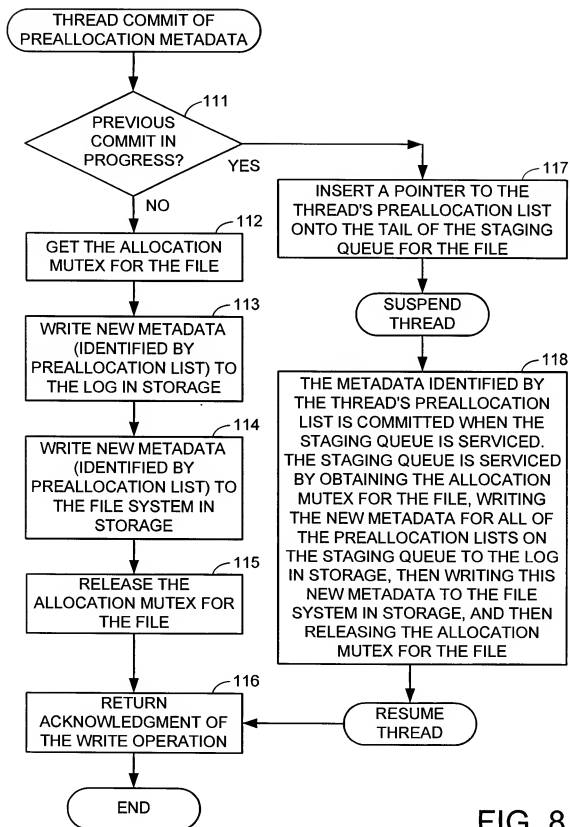


FIG. 8

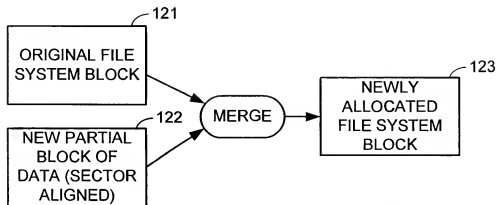


FIG. 9

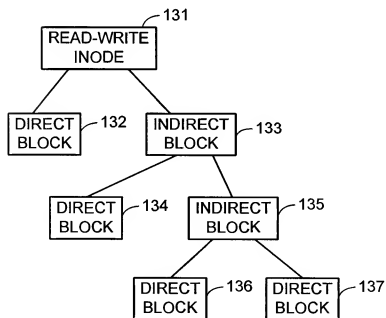


FIG. 10

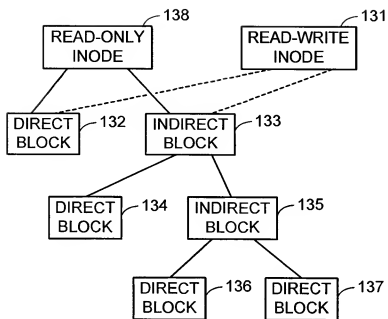


FIG. 11

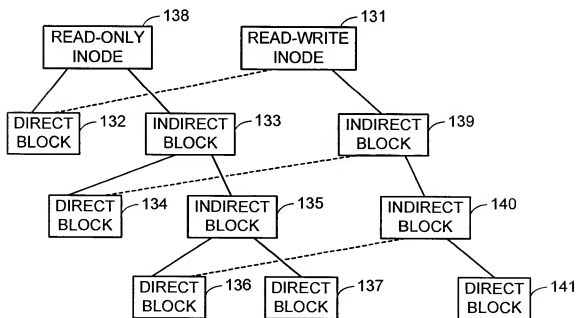


FIG. 12

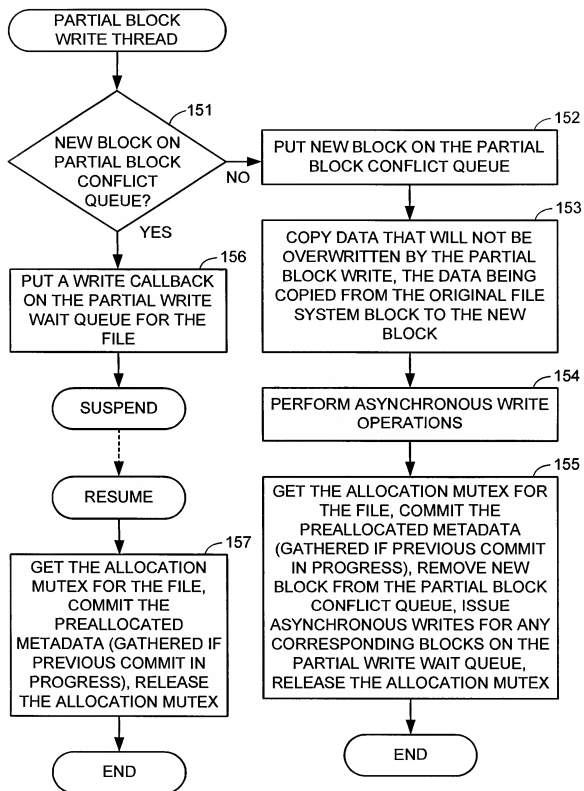


FIG. 13

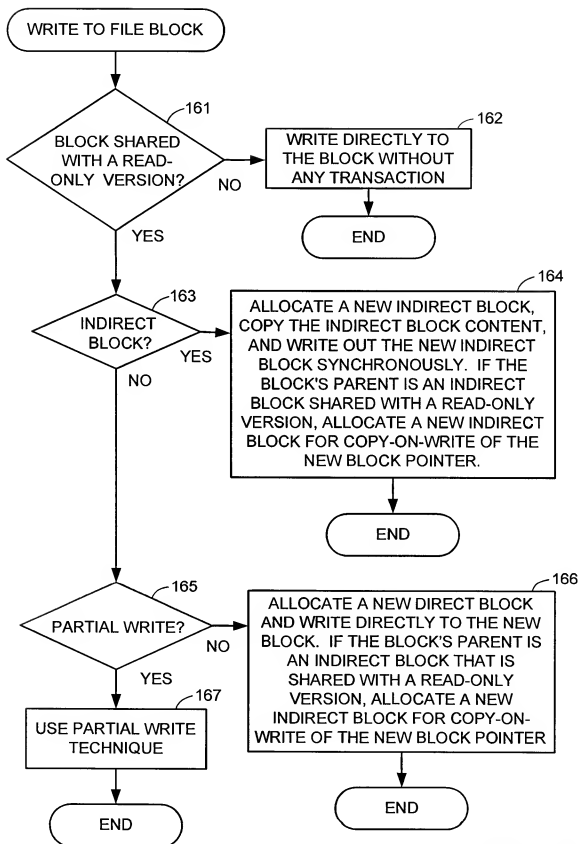


FIG. 14

EXHIBIT C

Auchterlonie, Richard

From: Auchterlonie, Richard
Sent: Thursday, September 04, 2003 12:59 PM
To: 'sfaibish@emc.com'
Cc: 'Clark_William@emc.com'; 'Mazzarella_Julie@emc.com'
Subject: EMC-03-062 A New Method for Enabling Simultaneous Parallel Writes to a Single File

Sorin:

Please find attached a second draft of the specification, drawings, declaration, and assignment for your patent application. Please review and pass it along to the other inventors as appropriate. If you would like any changes, please let me know and I will send you a revised draft. Otherwise, once all of the inventors have approved of the patent application, please contact Julie Mazzarella to arrange for execution of the patent application.

Thanks,

Richard C. Auchterlonie
Howrey Simon Arnold & White, LLP
750 Bering Drive
Houston, Texas 77057
Phone: (713) 787-1698
Fax: (713) 787-1440
AuchterlonieR@howrey.com



EMCR 100
Drawings.pdf



EMCR 100 PA
Parallel Writes.do...



EMCR100
Assignment.DOC



EMCR100
Declaration.DOC